

001 **AvatarPose: Avatar-guided 3D Pose Estimation**
002 **of Close Human Interaction**
003 **from Sparse Multi-view Videos**

004 Anonymous ECCV 2024 Submission

005 Paper ID #6979



006 **Fig. 1:** We propose AvatarPose, a method to estimate the 3D poses and shapes of multiple closely interacting people from multi-view videos. To this end, we first reconstruct the avatar of each individual and leverage the learned personalized avatars as priors to refine poses via color and silhouette rendering loss. We alternate between avatar refinement and pose optimization to obtain the final pose estimation.

007 **Abstract.** Despite progress in human motion capture, existing multi-view methods often face challenges in estimating the 3D pose and shape
008 of multiple closely interacting people. This difficulty arises from reliance
009 on accurate 2D joint estimations, which are hard to obtain due to occlu-
010 sions and body contact when people are in close interaction. To address
011 this, we propose a novel method leveraging the personalized implicit
012 neural avatar of each individual as a prior, which significantly improves
013 the robustness and precision of this challenging pose estimation task.
014 Concretely, the avatars are efficiently reconstructed via layered volume
015 rendering from sparse multi-view videos. The reconstructed avatar prior
016 allows for direct optimization of 3D poses based on color and silhouette
017 rendering loss, bypassing the issues associated with noisy 2D detections.
018 To handle interpenetration, we propose a collision loss on the overlap-
019 ing shape region of avatars to add penetration constraints. Moreover,
020 both 3D poses and avatars are optimized in an alternating manner. Our
021 experimental results demonstrate state-of-the-art performance on several
022 public datasets.

024 **Keywords:** human pose estimation · human close interaction · multi-
025 view pose estimation · avatar prior

026

1 Introduction

027 People frequently engage in interactions with each other in daily life to offer
 028 physical support or convey emotions. For AI systems to interpret 3D human
 029 interactions, a foundational step is to reconstruct 3D human poses and shapes
 030 using consumer-grade sensors like cameras. However, in closely interacting sce-
 031 narios, such as hugging or kissing, 3D pose estimations face challenges due to
 032 strong occlusions. To predict 3D poses from 2D images or videos, the problem
 033 of depth ambiguity gets worse when close contact is involved. Consequently, a
 034 multi-camera setup becomes essential to offer additional observations and to
 035 address depth ambiguity in pose estimation.

036 Despite the ubiquity of close human interactions, the study of estimating such
 037 human motions is scarce. Most previous multi-human methods with multi-view
 038 setup [2–4, 19, 22, 30, 57, 65, 71] focus on scenarios where people are at a distance
 039 from each other. Some methods [2, 4, 19, 22, 71] typically formulate the prob-
 040 lem as a cross-view matching problem, relying heavily on 2D joint estimations
 041 for subsequent 3D triangulation. These methods demonstrate high sensitivity to
 042 noisy or missing 2D joint estimations, particularly when occlusion occurs. An-
 043 other group of learning-based methods [57, 60, 62, 65] attempts to integrate 2D
 044 features from each view into a 3D voxel space and predict 3D human poses from
 045 identified 3D subvolumes of each individual. These methods are more robust to
 046 occlusion, but they encounter challenges with generalization and are sensitive to
 047 distribution shifts due to the lack of annotated 3D data. To correct abnormal or
 048 missing pose estimations, some methods [18, 21] leverage parametric body mod-
 049 els like SMPL [37] as full-body priors [18, 21] and fit these 3D models to 2D joint
 050 estimations. Although this method alleviates the issue of abnormal or missing
 051 joint estimations, it remains constrained by noisy 2D detections.

052 Embracing the challenging problem of pose estimation for close interactions,
 053 we propose a novel method to estimate the 3D poses and shapes of multiple
 054 people, observed from a sparse set of cameras. Our goal is to ensure that even
 055 in close contact, the estimated human poses and shapes are accurate and with-
 056 out interpenetrations. The key of our method is reconstructing implicit textured
 057 avatars of each individual in the scene and leveraging them as a strong person-
 058 alized prior for pose optimization (Figure 1). In contrast to relying solely on
 059 noisy 2D joint detections, this textured avatar prior enables us to leverage color
 060 and silhouette information for pose refinement, which significantly increases the
 061 robustness and accuracy of our method (Table 3). Meanwhile, the reconstructed
 062 avatar provides crucial geometric and appearance information to avoid collisions
 063 between individuals. Compared to methods using SMPL body shape [25, 42] to
 064 penalize collisions, our implicit avatar model contains a more detailed clothed
 065 shape and additional appearance cues and enables efficient computation of pen-
 066 etration loss. Due to the mutually beneficial relationship between avatar and
 067 pose, we alternate between pose optimization and avatar refinement.

068 More specifically, to accelerate the learning of avatars, we model each human
 069 individual in canonical space using an efficient neural radiance field variant in
 070 Instant NGP [43] and combine it with an efficient SMPL-based deformation

071 module [37]. To learn and render avatar models of multiple people, we adapt
 072 layered volume rendering [52, 69] to our avatar model. This adaptation allows
 073 us to jointly optimize all avatar models through a straightforward rendering
 074 loss. Once learned, the avatar can be animated and rendered based on pose
 075 parameters at interactive rates, thus naturally serving as an efficient personalized
 076 textured prior for pose optimization. With the learned personalized prior, we
 077 optimize pose via a novel objective function. Different from previous methods
 078 based on 2D reprojection error of joints, we directly optimize pose parameters via
 079 minimizing color and silhouette rendering losses while keeping the learned avatar
 080 model fixed. To prevent interpenetration between human individuals, a collision
 081 loss is introduced by penalizing the situation when a 3D point is occupied by
 082 multiple avatars. To remove artifacts in the initial avatar due to imperfect pose
 083 initialization, we further refine avatars based on optimized poses. Throughout
 084 the optimization process, both initial personalized avatar models and SMPL
 085 parameters are optimized in an alternating manner, motivated by the insight
 086 that accurate 3D pose estimations improve avatar learning, and improved avatar
 087 models, in turn, increase the precision of overall pose estimations.

088 We experimentally demonstrate that our method significantly outperforms
 089 previous state-of-the-art methods on several public datasets both quantitatively
 090 (Table 1 and Table 2) and qualitatively (Figure 3) especially when people are in
 091 close interaction. In summary, our contributions are:

- 092 – We propose a pipeline that efficiently creates implicit neural avatars of closely
 093 interacting people and leverages the learned avatars as priors to optimize
 094 poses.
- 095 – The avatar prior enables us to design a novel objective function that leverages
 096 color and mask rendering loss for pose optimization. We show the superiority
 097 of this loss function compared to the 2D reprojection error of 3D joints, which
 098 is used by most of the previous methods.
- 099 – Based on the learned avatar, a collision loss is introduced to avoid penetra-
 100 tion when individuals are in contact.

101 2 Related Work

102 2.1 Multi-Person 3D Pose Estimation

103 Despite significant progress in multi-human 2D pose estimation [9, 13, 15, 24, 34,
 104 48] and 3D pose estimation from monocular image or video [7, 31–33, 38–40, 45,
 105 55, 56, 58, 67, 68], the reconstruction accuracy is still limited due to depth ambi-
 106 guity and strong occlusions when humans are in close contact with each other.
 107 A multi-view setting [4, 5, 19, 22, 23, 35, 36, 53, 60, 62, 65, 71] helps to alleviate
 108 these challenges. One straightforward idea of most methods [2, 4, 19, 22, 71, 76]
 109 is to formulate the problem into cross-view matching and association problems.
 110 MVPose [19] performs 2D person parsing in each image and leverages cross-view
 111 person matching to infer 3D pose. 4DAssociation [71] additionally adds tracking

into this process to form a unified graph for associating 4D information. However, these methods are sensitive to the noisy estimation of 2D pose. In contrast to these matching-based methods, some recent methods [5, 14, 17, 23, 35, 49, 57, 60, 62, 65, 70, 75] directly learn deep neural networks to regress poses. Faster Voxelpose [65] employs the feature volume proposed by Voxelpose [57] and enhances computational efficiency. Graph [62] designs three graph neural network models for human center detection and pose estimation. MvP [60] simplifies the multi-person pose estimation by direct regression using the transformer model. A concurrent method CloseMocap [53] proposes to learn a model from a synthetic dataset simulating occlusion situations. However, heavily relying on the 2D or 3D features as input during training, these learning-based methods suffer from generalization issues when subjects, motions, and camera configurations change.

To further improve the robustness, a statistical parametric body model such as SMPL [37] is explored in [21] as a regularization prior for 3D joint refinement. Some follow-ups [18, 72] show that parametric models help in correcting implausible 3D pose estimates and filling in missing joints. However, this coarse body prior highly relies on aligning the 3D joints to 2D pose estimations, which are inaccurate when occlusion happens. Different from all of these previous methods, we explore the usage of personalized textured avatar models as priors to refine human poses. This prior enables us to leverage color and silhouette information from multi-view observation for refining poses of closely interacting humans.

2.2 3D Human Modeling

Parametric human body models [1, 28, 37, 44, 64] can represent minimally clothed human shapes by deforming a template mesh. It is challenging to extend this explicit representation for modeling clothed humans due to the fixed topology and resolution. To overcome this limitation, methods such as SNARF [12] and SCANimate [51] propose to model articulated human shapes based on 3D implicit representations. Many works [6, 8, 10, 20, 26, 29, 47, 59, 61] fit implicit neural fields to RGB or RGB-D videos by neural rendering to reconstruct the shape and appearance of a single human body. However, when applied to a multi-human scene, these methods are not able to achieve good fidelity due to strong occlusions. Recent methods including ST-NeRF [69] and [52] leverage layered neural representation to model multiple humans with sparse multi-view videos and thus can generate novel view synthesis of dynamic multiple humans. The main problem of all aforementioned approaches, however, is their high reliance on 3D human pose estimation [69, 73]: the deformation of the human model requires accurate human poses, which is hard to obtain when people are in close interaction. To address this challenge, our method is orthogonal to others, aiming at leveraging the learned avatar as priors for pose estimation in close interaction.

2.3 3D Human Datasets for Close Interaction

Most of the existing datasets [3, 11, 30] like Shelf and Panoptic studying multi-person pose estimation focus on the scene where people are at a distance from

each other and rarely involved in close interactions. To study close interactions among people, MultiHuman [73] dataset captures multi-person interaction with some close interactions and occlusions. ExPI [27] creates a multi-person extreme motion dataset with close interactions, but it focuses mainly on motion prediction for future frames instead of pose estimation from sparse views. CHI3D [25] captures two-person interaction datasets and proposes to learn a contact estimation module from annotations to improve the precision of pose estimation. The most recent work Hi4D [66] creates a challenging dataset of physically close human interaction and proposes a method to disentangle human bodies and estimate poses. However, this method relies on the 3D ground truth of clothed human meshes captured with expensive 3D body scanners. Due to the reliance on the limited annotated 3D data, this method faces challenges in generalization with different people and camera configurations. In contrast to all the works before, we intend to solve the problem of pose estimation in close human interaction without requiring accurate 3D scans or other training data.

3 Method

Given a dynamic scene captured by a sparse set of RGB cameras, our goal is to estimate the 3D pose and shape of multiple people even if they interact closely. To address this challenging task, our key idea is to first reconstruct the personalized avatar of each individual in the scene and leverage them as a strong prior to refine the appearance and pose in an alternating manner. An overview of our method is shown in Figure 2.

We first introduce an efficient pipeline to create avatars of multiple people in a scene (Section 3.1 and Figure 2(a)). Specifically, we leverage an accelerated neural radiance field to represent the shape and appearance of each individual in canonical space and deform it at an interactive rate. We then adapt layered volume rendering to our pipeline, which composites the rendering of avatars into one image, thus enabling direct learning from multi-view video inputs.

Thanks to the learned avatar prior of each individual, we can enhance 3D pose optimization via a combination of RGB and silhouette rendering loss (Section 3.2 and Figure 2(b)). While previous work highly relies on noisy 2D joint detection, we show that employing such color and silhouette information can largely increase precision and robustness. Moreover, a collision loss is introduced to avoid interpenetration. Finally, we alternate between avatar learning and pose optimization to generate complete and accurate 3D human poses.

3.1 Multi-Avatar Prior Learning

Avatar Model We represent each human individual in canonical space using an accelerated neural radiance field [43] and model shape-aware articulated deformation based on SMPL [37].

- **Canonical Appearance Representation:** To model human shape and appearance, we create canonical radiance field $\bar{\mathbf{F}}_{\sigma_f}^{(l)}$ for each human instance

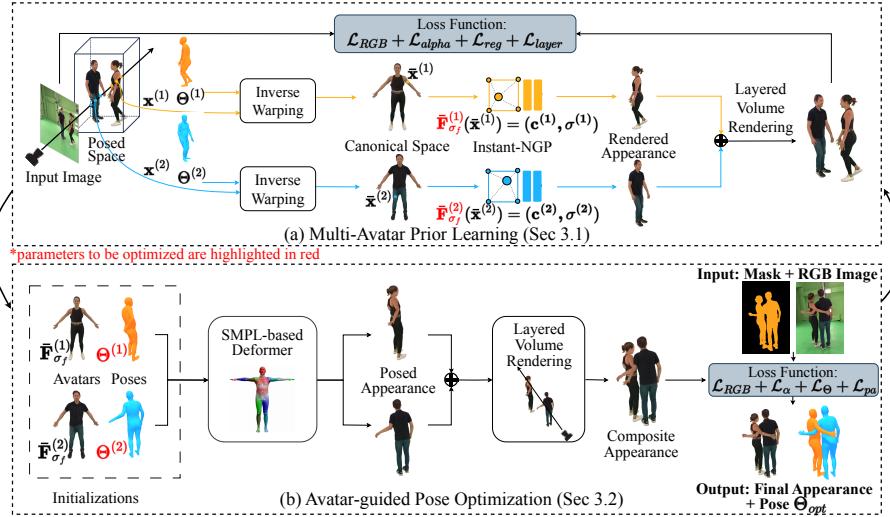


Fig. 2: Method Overview: Our method consists of two modules: (a) *Multi-Avatar Prior Learning*: Given the input multi-view images and estimated poses $\Theta^{(l)}$, we sample points $\mathbf{x}^{(l)}$ for each individual l along the rays in posed space and warp these points into canonical space and calculate their color $\mathbf{c}^{(l)}$ and density $\sigma^{(l)}$ via the canonical appearance network $\bar{\mathbf{F}}_{\sigma_f}^{(l)}$. We leverage layered volume rendering [69] to attain the final pixel color and compare it with the original input image to optimize the parameters of avatars. (b) *Avatar-guided Pose Optimization*: Given learned avatar model $\bar{\mathbf{F}}_{\sigma_f}^{(l)}$ and initial poses $\Theta^{(l)}$ of each individual l , we deform all of the avatars based on SMPL-based deformer and render them jointly via layered volume rendering. We compare the composite rendering with input observation and minimize the RGB and mask rendering loss to optimize poses. A collision loss is additionally introduced to avoid interpenetration. Finally, we alternate between two modules to obtain the final result. For clarity, the parameters to be optimized are marked as red in each module.

195 $l \in [1, L]$, where L is the number of humans in the scene. $\bar{\mathbf{F}}_{\sigma_f}^{(l)}$ takes a 3D
 196 point $\bar{\mathbf{x}}^{(l)}$ as input and predicts its density $\sigma^{(l)}$ and color $\mathbf{c}^{(l)}$. Following
 197 Instant-NGP [43] and InstantAvatar [29] to accelerate the rendering, $\bar{\mathbf{F}}_{\sigma_f}^{(l)}$ is
 198 parameterized via using a hash table to store feature grids at different scales.

- 199 – **Pose Representation:** We represent the 3D pose and underlying body
 200 shape for all human instances by SMPL parameters $\Theta = \{\Theta^{(l)}\}_{l \in [1, L]}$. For
 201 each human l , $\Theta^{(l)} = \{\beta^{(l)}, \theta^{(l)}, t^{(l)}\}$ contains shape parameters $\beta^{(l)} \in \mathbb{R}^{10}$,
 202 pose parameters $\theta^{(l)} \in \mathbb{R}^{72}$ and translation $t^{(l)} \in \mathbb{R}^3$ of SMPL.
- 203 – **Deformer:** To enable animation given targeted poses $\Theta^{(l)}$, we require the
 204 radiance field in the posed space. Given a point $\mathbf{x}^{(l)}$ in deformed space of hu-
 205 man l , we determine the corresponding canonical point $\bar{\mathbf{x}}^{(l)}$ by inverse linear
 206 blend skinning(LBS) [37]: $\bar{\mathbf{x}}^{(l)}(\mathbf{x}^{(l)}, \Theta^{(l)}) = (\sum_{i=1}^{n_b} w_i(\Theta^{(l)}) \mathbf{B}_i(\Theta^{(l)}))^{-1} \mathbf{x}^{(l)}$,
 207 where \mathbf{B}_i is the rigid bone transformation matrix for joint $i \in \{1, \dots, n_b\}$
 208 under pose $\Theta^{(l)}$. w_i is the skinning weights of the nearest neighbor of $\mathbf{x}^{(l)}$ in
 209 the deformed SMPL vertices. We obtain the radiance field at the point $\mathbf{x}^{(l)}$

210 by evaluating the canonical appearance field at the corresponding point $\bar{\mathbf{x}}^{(l)}$.
 211

212 *Layered Volume Rendering* To obtain the pixel value for a ray $\mathbf{r} \in \mathcal{R}$, we raycast
 213 every human instance separately with a layered rendering strategy similar to
 214 ST-NeRF [69]. Specifically, we first calculate the intersection points between the
 215 ray and the 3D bounding box of each human instance and uniformly sample
 216 points in each bounding box. To distinguish different identities, we assign each
 217 sampled point \mathbf{x}_i a one-hot representation $\mathbf{m}_i = [m_i^{(1)}, \dots, m_i^{(L)}]$ to indicate
 218 which human identity it belongs to. After sorting all sampled points by their
 219 depth values and calculating their corresponding color \mathbf{c}_i and density σ_i from
 220 the avatar model, if $m_i^{(l)} = 1$, we compute

$$221 \quad \mathbf{c}_i, \sigma_i = \bar{\mathbf{F}}_{\sigma_f}^{(l)}(\bar{\mathbf{x}}_i(\mathbf{x}_i, \Theta^{(l)})). \quad (1) \quad 221$$

223 The color of each ray is computed via numerical integration [41].
 224

$$224 \quad \hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j < i} (1 - \alpha_j) \mathbf{c}_i \quad \alpha_i = 1 - \exp(-\sigma_i \delta_i), \quad (2) \quad 224$$

225 where δ_i is the distance between samples. The accumulated alpha value,
 226 which represents ray opacity, can be computed via:
 227

$$228 \quad \alpha(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j < i} (1 - \alpha_j). \quad (3) \quad 228$$

229 For each human identity $l \in [1, L]$, the corresponding instance ray opacity
 230 can be calculated via:
 231

$$232 \quad \alpha^{(l)}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j < i} (1 - \alpha_j) m_i^{(l)}. \quad (4) \quad 232$$

233 *Training* The overall training process is shown in Figure 2(a). For training avatar
 234 layers, we minimize the Huber loss ρ between the predicted pixel color $\hat{\mathbf{C}}(\mathbf{r})$ and
 235 the ground truth pixel color $\mathbf{C}_{gt}(\mathbf{r})$:
 236

$$237 \quad \mathcal{L}_{RGB} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \rho(\|\hat{\mathbf{C}}(\mathbf{r}) - \mathbf{C}_{gt}(\mathbf{r})\|). \quad (5) \quad 237$$

238 Since instance segmentation of human performers is hard to obtain and is not
 239 accurate, we choose foreground segmentation as our mask supervision, which is
 240 obtained via SAM-Track [16]. We apply a loss for optimizing the rendered alpha
 241 values α to reduce the artifacts in the floating area:
 242

$$243 \quad \mathcal{L}_{alpha} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} (\alpha(\mathbf{r}) - \alpha_{SAM}(\mathbf{r}))^2. \quad (6) \quad 243$$

Following [52], we add a regularization loss for instance alpha values to make sure every pixel can only be rendered from one human layer:

$$\mathcal{L}_{layer} = -\frac{1}{L |\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \sum_{l=1}^L \alpha^{(l)}(\mathbf{r}) \log(\alpha^{(l)}(\mathbf{r})). \quad (7)$$

Similar to [29], we also add hard surface and density regularization terms in the learning process. More training details can be found in the Supp Mat.

3.2 Avatar-guided Pose Optimization

Equipped with learned avatars obtained in Section 3.1, we aim to estimate the 3D shape and pose of multiple humans with close physical contact. To achieve this, we leverage the avatars as priors to handle the challenging pose ambiguities caused by contact. Specifically, we first initialize pose parameters using an off-the-shelf 3D pose estimator [71] and refine the pose via a rendering loss between rendered posed avatars and 2D observations (Figure 2(b)). Since the initial imperfect pose estimations may cause artifacts in avatar reconstruction, we further refine the weights of the avatar model using the optimized pose. This process is formulated as an alternating optimization to refine both poses and avatars.

Initialization The initial 3D pose proposals are estimated by leveraging the off-the-shelf 3D human pose estimator [71]. After that, we register a SMPL model to the estimated 3D joints to obtain the initial pose parameters Θ_0 . Given these estimated poses, the initial avatar model of each individual is further learned from multi-view videos.

Objective To leverage avatars as priors to tackle challenges caused by contact, we optimize the following objective:

$$\begin{aligned} \mathcal{L}(\Theta) = & \lambda_{RGB} \mathcal{L}_{RGB}(\Theta) + \lambda_\alpha \mathcal{L}_\alpha(\Theta) \\ & + \lambda_{reg} \mathcal{L}_{reg}(\Theta) + \lambda_{pa} \mathcal{L}_{pa}(\Theta). \end{aligned} \quad (8)$$

Here, Θ is the SMPL parameters to be optimized for all avatars, which is also consistent with the pose parameters to deform the avatar model (Section 3.1). Following [61], we represent Θ as $\Theta_0 + MLP(\Theta_0)$ and refine poses by changing the parameters of the neural network. This representation empirically shows more robust results compared to directly optimizing SMPL parameters.

Different from previous methods [18, 21, 57, 71], the personalized multi-avatar prior allows us to leverage appearance and silhouette information to refine initial poses. Specifically, we calculate \mathcal{L}_{RGB} to ensure the color consistency between the rendered pixel $\hat{\mathbf{C}}(\mathbf{r}, \Theta)$ (Equation (1) and Equation (2)) of deformed avatar with poses Θ and the corresponding ground-truth pixel color $\mathbf{C}_{gt}(\mathbf{r})$.

$$\mathcal{L}_{RGB}(\Theta) = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \rho(\|\hat{\mathbf{C}}(\mathbf{r}, \Theta) - \mathbf{C}_{gt}(\mathbf{r})\|). \quad (9)$$

282 Additionally, a cross-entropy loss \mathcal{L}_α is introduced to ensure that the rendered
 283 mask $\alpha(\mathbf{r}, \Theta)$ (Equation (3)) of the reposed avatar is aligned with the estimated
 284 SAM-Track mask $\alpha_{SAM}(\mathbf{r})$ by:

$$285 \quad \mathcal{L}_\alpha(\Theta) = - \sum_{\mathbf{r} \in \mathcal{R}} \alpha_{SAM}(\mathbf{r}) \log(\alpha(\mathbf{r}, \Theta)) \quad (10) \quad 285$$

286 To penalize the unnatural poses and avoid elbows and knees bending in the
 288 wrong direction, we add an L2 regularization term and combine it with the pose
 289 prior in SMPLify [7] to constrain physically implausible joint rotation:

$$290 \quad \mathcal{L}_{reg}(\Theta) = \|\Theta\|_2 + \lambda \sum_{i \in I} \exp(\Theta_i), \quad (11) \quad 290$$

291 where I is the set of pose indices corresponding to elbows and knees.

293 A key challenge to correctly estimate poses in close interaction is to handle
 294 interpenetration. Since every avatar is modeled separately, the surfaces tend to
 295 intersect when they are in contact. To handle this, we first select sampled points
 296 inside multiple instances as $\mathcal{S} = \{\mathbf{x}_i \mid \alpha_i^{(p)} > 0, \alpha_i^{(q)} > 0, p, q \in [1, L], p \neq q\}$ ($\alpha_i^{(p)}$
 297 is calculated from Equation (1) and Equation (2) corresponding with a point \mathbf{x}_i
 298 with $m_i^{(p)} = 1$). We then propose a collision loss \mathcal{L}_{pa} for penalizing penetration:

$$299 \quad \mathcal{L}_{pa}(\Theta) = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x}_i \in \mathcal{S}} \alpha_i^{(p)}(\Theta) \alpha_i^{(q)}(\Theta) \quad (12) \quad 299$$

300 Intuitively, this loss guarantees every sample point in 3D space can not be
 302 occupied by multiple avatar models simultaneously, which guarantees better pose
 303 estimation in close contact.

304 *Alternating Optimization* Since artifacts sometimes appear on initial avatars due
 305 to imperfect pose initialization, we further refine avatars based on optimized
 306 poses via minimizing the loss function in Section 3.1. Finally, the optimization
 307 of poses and avatars is formulated in an alternating fashion for N steps. More
 308 details of optimization can be found in the Supp Mat.

309 4 Experiments

310 *Datasets.* We mainly evaluate our proposed method on Hi4D [66] and CHI3D [25]
 311 Dataset, which are challenging datasets of two humans in close interaction. To
 312 further demonstrate the generalization ability of our method for more than two
 313 people, we also evaluate our method on Shelf [3] and MultiHuman [52] Dataset,
 314 which includes three or four people. More details are shown in the Supp Mat.

315 *Metrics.* We use the Mean Per Joint Position Error (MPJPE) [57] to measure
 316 the distance between the ground truth 3D poses and the estimated poses. We
 317 also choose the Percentage of Correct Parts (PCP3D) metric [19] to calculate
 318 the percentage of correct parts. AP_K [57] and Recall [21] are also leveraged to
 319 evaluate performance. More details are shown in the Supp Mat.

Baselines. We compare our method on the Hi4D dataset with state-of-the-art methods in three categories discussed in Section 2.1. For learning-based methods, we choose Graph [63], MvP [60], and Faster VoxelPose [65] and fine-tune these models on a subset of the evaluating datasets. For pure association-based methods, we choose 4DAssociation [71] and MVPose* [19]. Based on initial results, MVPose [18] added temporal tracking and SMPL prior to MVPose* [19], which is regarded as an SMPL-guided method. More details about baselines can be found in the Supp Mat.

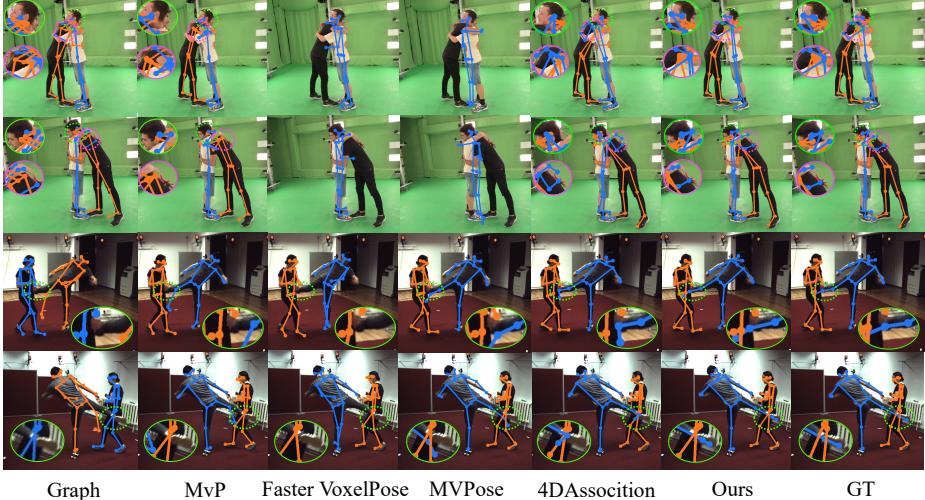


Fig. 3: Qualitative Comparison with SotA methods [19, 60, 63, 65, 71] on Hi4D and CHI3D. We show two examples from the Hi4D and CHI3D datasets compared with Graph, MvP, Faster VoxelPose, MVPose, and 4DAssociation. For each example, we show 2D projections on two sampled views.

4.1 Comparison to SotA

Table 1 and Table 2 summarizes our quantitative comparisons on Hi4D and CHI3D with SotA (State-of-the-art) methods. Our method largely outperforms the other SotA methods in all of the metrics including MPJPE, PCP3D, and AP_K . More comparisons on Shelf and MultiHuman are shown in the Supp Mat.

Comparison with Graph [62], MvP [60], Faster VoxelPose [65]. Our method achieves much better quantitative results compared to Graph, MvP, and Faster VoxelPose on Hi4D and CHI3D. These learning-based methods are prone to overfitting to training pose distribution, resulting in an inability to accurately estimate the challenging poses of interacting actors, such as touching shoulders or legs. More specifically, Faster VoxelPose sometimes causes missing actors in close contact, leading to a relatively low recall. Graph and MvP perform better

| Method | MPJPE(mm) ↓ | PCP(%) ↑ | AP ₅₀ ↑ | AP ₁₀₀ ↑ | Recall(%) ↑ |
|-----------------------|--------------|--------------|--------------------|---------------------|--------------|
| MvP [60] | 92.77 | 74.14 | 41.59 | 63.86 | 93.84 |
| Graph [62] | 89.62 | 71.55 | 44.75 | 67.33 | 93.31 |
| Faster VoxelPose [65] | 68.40 | 73.67 | 44.05 | 68.70 | 83.55 |
| MVPose* [19] | 53.05 | 87.57 | 67.97 | 80.28 | 93.80 |
| MVPose [18] | 42.63 | 90.76 | 71.79 | 90.19 | 93.30 |
| 4DAssociation [71] | 41.29 | 88.62 | 80.87 | 97.27 | 98.78 |
| Ours | 32.10 | 96.90 | 91.48 | 97.33 | 98.78 |

Table 1: Quantitative Comparison with SotA on the Hi4D [66] Dataset (8 views). We compare our method with MvP [60], Graph [62], Faster VoxelPose [65], MVPose* [19], MVPose [18] and 4DAssociation [71]. We report MPJPE, PCP, AP_K, and Recall metric for all methods.

| Method | MPJPE(mm) ↓ | PCP(%) ↑ | AP ₅₀ ↑ | AP ₁₀₀ ↑ | Recall(%) ↑ |
|-----------------------|--------------|--------------|--------------------|---------------------|--------------|
| MvP [60] | 55.38 | 89.47 | 63.58 | 92.53 | 99.06 |
| Graph [62] | 45.33 | 92.46 | 74.02 | 95.25 | 99.17 |
| Faster VoxelPose [65] | 67.81 | 78.41 | 29.28 | 82.88 | 93.34 |
| MVPose* [19] | 50.42 | 90.39 | 69.13 | 75.72 | 88.72 |
| MVPose [18] | 34.05 | 93.35 | 79.94 | 86.91 | 88.18 |
| 4DAssociation [71] | 37.47 | 99.30 | 89.66 | 98.67 | 99.85 |
| Ours | 32.98 | 99.79 | 93.20 | 99.79 | 99.85 |

Table 2: Quantitative Comparison with SotA on the CHI3D [25] Dataset (4 views). We compare our method with Faster VoxelPose [65], MVPose* [19], MVPose [18] and 4DAssociation [71]. We report MPJPE, PCP, AP_K, and Recall metric for all methods.

in recall, but Graph fails to consistently track the actors across frames, and MvP results in many misaligned joints between actors.

Comparison with MVPose [18]. Compared with MVPose, the MPJPE and precision of our method are much better. Specifically, when the threshold of precision becomes smaller, the gap of precision becomes larger. This is because the precision of fitting the SMPL body heavily relies on 2D joint detections, which are noisy and inaccurate when close interaction happens. In contrast, our method takes full advantage of color and silhouette rendering loss to optimize poses, leading to robustness to occlusions. Furthermore, this top-down method also cannot detect closely interacting actors correctly due to strong occlusions shown in Figure 3.

Comparison with 4DAssociation [71]. Finally, we compare our method with the bottom-up association method. As shown in Table 1 and Table 2, our method outperforms 4DAssociation in most metrics. Figure 3 shows that when actors are close, this bottom-up method is inclined to associate joints with the wrong human instances. This is because they solve the joint association with a greedy algorithm, which is sensitive to missing and inaccurate 2D joint detections.

357 4.2 Additional Qualitative Samples

358 Figure 4 demonstrates more qualitative results of our method on Hi4D, CHI3D,
 359 and MultiHuman Real-Cap with challenging and close interactions among 2 or
 360 3 people. we also demonstrate results on the Shelf [3] dataset, which contains 4
 361 people without close contact. More results can be found in the Supp Mat.

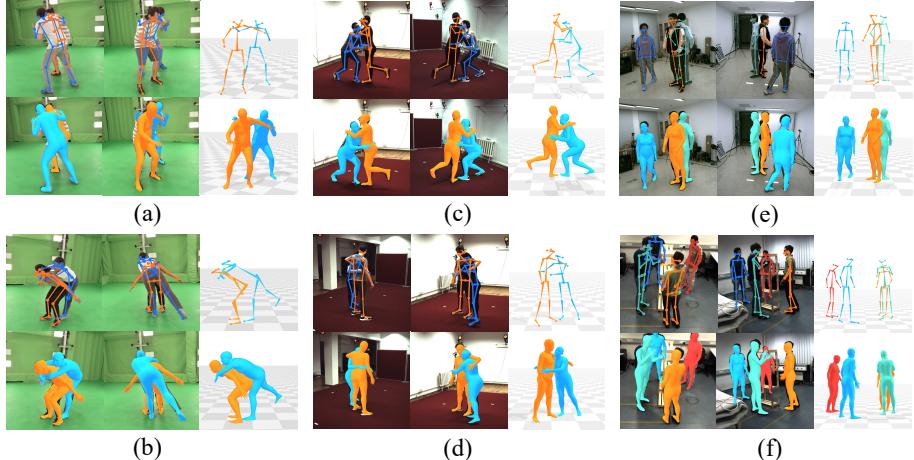


Fig. 4: Qualitative Results of our method on Hi4D (a)(b), CHI3D (c)(d), MultiHuman Real-Cap (e), and Shelf (f). The left and middle columns in each sub-figure show the 2D projections of the estimated 3D skeletons and SMPL body meshes on two views. The right column in each sub-figure demonstrates skeletons and SMPL bodies in 3D scenes.

362 4.3 Ablation Study

363 To validate the effectiveness of our method, we conduct a detailed analysis of
 364 different design choices of our algorithm. All the experiments are conducted on
 365 the Hi4D Dataset.

366 *Comparison with SMPL Body Prior.* To validate the effectiveness of our person-
 367 alized avatar prior, we compare our method with a baseline that optimizes SMPL
 368 parameters to align reprojections of 3D joints to 2D observations. As shown in
 369 Table 3, our method outperforms the SMPL prior baseline by a significant margin.
 370 This is due to fitting errors when 2D pose estimation is not accurate. In
 371 contrast, thanks to the learned avatar prior, our method can leverage color and
 372 silhouette rendering loss to refine poses, without high reliance on joint detection
 373 and alignment which are not accurate when occlusion happens. Figure 5 shows
 374 a qualitative comparison. We observe that the pose of the arm in the baseline
 375 method is wrongly estimated and even leads to penetration, whereas our method
 376 still reconstructs the poses correctly.

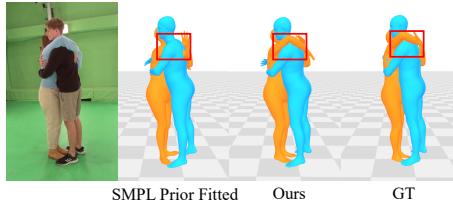


Fig. 5: Comparison with SMPL Body Prior. Only fitting SMPL to 2D observations, some joints in close contact such as arms are incorrectly estimated and even cause intersections between body surfaces. In contrast, our personalized prior enables accurate estimation of poses.

Color and Mask Loss. To demonstrate the effectiveness of the color and silhouette rendering loss in our optimization process, we design baselines without RGB loss and mask loss respectively for comparison. Table 3 shows that the RGB loss significantly improves pose optimization, and the optimization process will completely deviate from the correct trajectory in the absence of RGB loss. Mask loss is proven to slightly increase the accuracy by adding additional constraints on the rendered human silhouette.

| Method | MPJPE(mm) ↓ | PCP(%) ↑ | AP ₅₀ ↑ | AP ₁₀₀ ↑ |
|---------------------------|--------------|--------------|--------------------|---------------------|
| Ours (SMPL fitted) | 40.41 | 95.10 | 84.04 | 94.37 |
| Ours (w/o RGB loss) | 78.40 | 83.26 | 17.84 | 74.55 |
| Ours (w/o Mask loss) | 31.00 | 97.33 | 90.28 | 99.06 |
| Ours (Joint Optimization) | 66.04 | 84.56 | 22.55 | 76.70 |
| Ours | 29.37 | 98.02 | 96.79 | 99.06 |

Table 3: Quantitative Ablation Results. Ablations to evaluate our method with only the SMPL fitted method, our method without RGB loss and without Silhouette loss, and our method without alternating optimization.

Alternating Optimization. To show the advantage of alternating optimization, a joint optimization is selected for comparison, which is widely used in avatar reconstruction [20, 26]. In Table 3, we observe that our method largely outperforms the baseline. As shown by the qualitative results in Figure 6, the final avatar of our ablated baselines suffers from artifacts and floating points around contact body parts. This imperfect avatar in turn causes the wrong 3D pose estimations. In contrast, our method can faithfully reconstruct avatars and poses under challenging poses.

Penetration Loss. Our penetration loss serves an important role in avoiding interpenetration. Comparing our method to the ablated version where we remove this loss, Figure 7 shows that the SMPL body of one person partially intersects with the other person in the contact area. This is also confirmed by the rendering result of the avatar. By penalizing the collision of density fields of avatars in

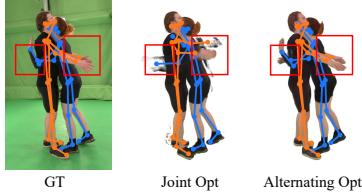


Fig. 6: Ablation of Alternating Optimization. We show the results of rendered avatars and projections of the estimated 3D poses. Joint optimization suffers from artifacts around the contact part and in turn causes wrong pose estimations. In contrast, ours reconstructs both avatars and poses correctly.

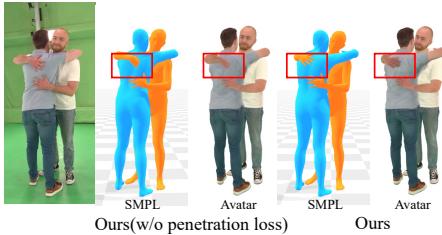


Fig. 7: Ablation of Penetration Loss. Without penetration loss, both the avatar and underlying SMPL body tend to have collisions on surfaces.

3D space, our method largely reduces the penetration and thus achieves more accurate pose estimation.

399 4.4 Limitations and Future Work

400 We do not model hands in our avatar model, it will be a promising direction
 401 to integrate hand models [46, 50] into the personalized avatar. The current optimi-
 402 zation pipeline is not fast enough, which can be accelerated with more efficient
 403 representation [74] or more powerful optimization tools [54]. We believe it will
 404 be interesting future work to adapt our pipeline to estimate 3D poses of closely
 405 interacting people from in-the-wild and monocular videos and images. More dis-
 406 cussions about limitations and future work can be found in the Supp Mat.

407 5 Conclusion

408 In this paper, we propose AvatarPose, a novel method to estimate the 3D poses
 409 of multiple people in close interaction from sparse multi-view videos. Unlike pre-
 410 vious methods leveraging SMPL body prior, we first reconstruct the avatar of
 411 each individual and leverage the avatars as personalized priors to guide pose
 412 optimization. The avatar prior enables us to use color and silhouette observa-
 413 tions, instead of relying on noisy 2D joint detections. A collision loss is also
 414 introduced to constrain penetration in close contact. Our method outperforms
 415 SotA methods significantly on public datasets of close human interactions. The
 416 code associated with this paper will be released upon acceptance.

417

References

417

- 418 1. Anguelov, D., Srinivasan, P., Koller, D., Thrun, S., Rodgers, J., Davis, J.: Scape:
419 shape completion and animation of people. In: ACM SIGGRAPH 2005 Papers, pp.
420 408–416 (2005) 4
- 421 2. Belagiannis, V., Amin, S., Andriluka, M., Schiele, B., Navab, N., Ilic, S.: 3d pic-
422 torial structures for multiple human pose estimation. In: Proceedings of the IEEE
423 conference on computer vision and pattern recognition. pp. 1669–1676 (2014) 2, 3
- 424 3. Belagiannis, V., Amin, S., Andriluka, M., Schiele, B., Navab, N., Ilic, S.: 3d pic-
425 torial structures for multiple human pose estimation. In: Proceedings of the IEEE
426 conference on computer vision and pattern recognition. pp. 1669–1676 (2014) 2,
427 4, 9, 12
- 428 4. Belagiannis, V., Amin, S., Andriluka, M., Schiele, B., Navab, N., Ilic, S.: 3d pic-
429 torial structures revisited: Multiple human pose estimation. IEEE transactions on
430 pattern analysis and machine intelligence **38**(10), 1929–1942 (2015) 2, 3
- 431 5. Benzine, A., Chabot, F., Luvison, B., Pham, Q.C., Achard, C.: Pandanet:
432 Anchor-based single-shot multi-person 3d pose estimation. In: Proceedings of the
433 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6856–
434 6865 (2020) 3, 4
- 435 6. Bhatnagar, B.L., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Combining im-
436 plicit function learning and parametric models for 3d human reconstruction. In:
437 Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August
438 23–28, 2020, Proceedings, Part II **16**. pp. 311–329. Springer (2020) 4
- 439 7. Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep
440 it smpl: Automatic estimation of 3d human pose and shape from a single image.
441 In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The
442 Netherlands, October 11–14, 2016, Proceedings, Part V **14**. pp. 561–578. Springer
443 (2016) 3, 9
- 444 8. Burov, A., Nießner, M., Thies, J.: Dynamic surface function networks for clothed
445 human bodies. In: Proceedings of the IEEE/CVF International Conference on
446 Computer Vision. pp. 10754–10764 (2021) 4
- 447 9. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estima-
448 tion using part affinity fields. In: Proceedings of the IEEE conference on computer vision
449 and pattern recognition. pp. 7291–7299 (2017) 3
- 450 10. Chen, J., Zhang, Y., Kang, D., Zhe, X., Bao, L., Jia, X., Lu, H.: Animatable neural
451 radiance fields from monocular rgb videos. arXiv preprint arXiv:2106.13629 (2021)
452 4
- 453 11. Chen, L., Ai, H., Chen, R., Zhuang, Z., Liu, S.: Cross-view tracking for multi-human
454 3d pose estimation at over 100 fps. In: Proceedings of the IEEE/CVF conference
455 on computer vision and pattern recognition. pp. 3279–3288 (2020) 4
- 456 12. Chen, X., Zheng, Y., Black, M.J., Hilliges, O., Geiger, A.: Snarf: Differentiable
457 forward skinning for animating non-rigid neural implicit shapes. In: Proceedings
458 of the IEEE/CVF International Conference on Computer Vision. pp. 11594–11604
459 (2021) 4
- 460 13. Chen, Y., Wang, Z., Peng, Y., Zhang, Z., Yu, G., Sun, J.: Cascaded pyramid
461 network for multi-person pose estimation. In: Proceedings of the IEEE conference
462 on computer vision and pattern recognition. pp. 7103–7112 (2018) 3
- 463 14. Chen, Y., Gu, R., Huang, O., Jia, G.: Vtp: volumetric transformer for multi-view
464 multi-person 3d pose estimation. Applied Intelligence pp. 1–12 (2023) 4

- 465 15. Cheng, B., Xiao, B., Wang, J., Shi, H., Huang, T.S., Zhang, L.: Higherhrnet: Scale-
466 aware representation learning for bottom-up human pose estimation. In: Proceedings
467 of the IEEE/CVF conference on computer vision and pattern recognition. pp.
468 5386–5395 (2020) 3 465
- 469 16. Cheng, Y., Li, L., Xu, Y., Li, X., Yang, Z., Wang, W., Yang, Y.: Segment and
470 track anything. arXiv preprint arXiv:2305.06558 (2023) 7 466
- 471 17. Choudhury, R., Kitani, K.M., Jeni, L.A.: Tempo: Efficient multi-view pose esti-
472 mation, tracking, and forecasting. In: Proceedings of the IEEE/CVF International
473 Conference on Computer Vision. pp. 14750–14760 (2023) 4 467
- 474 18. Dong, J., Fang, Q., Jiang, W., Yang, Y., Bao, H., Zhou, X.: Fast and robust multi-
475 person 3d pose estimation and tracking from multiple views. In: T-PAMI (2021)
476 2, 4, 8, 10, 11 468
- 477 19. Dong, J., Jiang, W., Huang, Q., Bao, H., Zhou, X.: Fast and robust multi-person 3d
478 pose estimation from multiple views. In: Proceedings of the IEEE/CVF conference
479 on computer vision and pattern recognition. pp. 7792–7801 (2019) 2, 3, 9, 10, 11 469
- 480 20. Dong, Z., Guo, C., Song, J., Chen, X., Geiger, A., Hilliges, O.: Pina: Learning a
481 personalized implicit neural avatar from a single rgb-d video sequence. In: Proceed-
482 ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
483 pp. 20470–20480 (2022) 4, 13 470
- 484 21. Dong, Z., Song, J., Chen, X., Guo, C., Hilliges, O.: Shape-aware multi-person pose
485 estimation from multi-view images. In: Proceedings of the IEEE/CVF Interna-
486 tional Conference on Computer Vision. pp. 11158–11168 (2021) 2, 4, 8, 9 471
- 487 22. Ershadi-Nasab, S., Noury, E., Kasaei, S., Sanaei, E.: Multiple human 3d pose es-
488 timation from multiview images. Multimedia Tools and Applications 77, 15573–
489 15601 (2018) 2, 3 472
- 490 23. Fabbri, M., Lanzi, F., Calderara, S., Alletto, S., Cucchiara, R.: Compressed vol-
491 umetric heatmaps for multi-person 3d pose estimation. In: Proceedings of the
492 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7204–
493 7213 (2020) 3, 4 473
- 494 24. Fang, H.S., Xie, S., Tai, Y.W., Lu, C.: Rmpe: Regional multi-person pose estima-
495 tion. In: Proceedings of the IEEE international conference on computer vision. pp.
496 2334–2343 (2017) 3 474
- 497 25. Fieraru, M., Zanfir, M., Oneata, E., Popa, A.I., Olaru, V., Sminchisescu, C.: Three-
498 dimensional reconstruction of human interactions. In: Proceedings of the
499 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7214–
500 7223 (2020) 2, 5, 9, 11 475
- 501 26. Guo, C., Jiang, T., Chen, X., Song, J., Hilliges, O.: Vid2avatar: 3d avatar re-
502 construction from videos in the wild via self-supervised scene decomposition. In:
503 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
504 nition. pp. 12858–12868 (2023) 4, 13 476
- 505 27. Guo, W., Bie, X., Alameda-Pineda, X., Moreno-Noguer, F.: Multi-person extreme
506 motion prediction. In: Proceedings of the IEEE/CVF Conference on Computer
507 Vision and Pattern Recognition. pp. 13053–13064 (2022) 5 477
- 508 28. Hasler, N., Stoll, C., Sunkel, M., Rosenhahn, B., Seidel, H.P.: A statistical model of
509 human pose and body shape. In: Computer graphics forum. vol. 28, pp. 337–346.
510 Wiley Online Library (2009) 4 478
- 511 29. Jiang, T., Chen, X., Song, J., Hilliges, O.: Instantavatator: Learning avatars from
512 monocular video in 60 seconds. In: Proceedings of the IEEE/CVF Conference on
513 Computer Vision and Pattern Recognition. pp. 16922–16932 (2023) 4, 6, 8 479

- 514 30. Joo, H., Liu, H., Tan, L., Gui, L., Nabbe, B., Matthews, I., Kanade, T., Nobuhara,
515 S., Sheikh, Y.: Panoptic studio: A massively multiview system for social motion
516 capture. In: Proceedings of the IEEE International Conference on Computer Vi-
517 sion. pp. 3334–3342 (2015) 2, 4
- 518 31. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human
519 shape and pose. In: Proceedings of the IEEE conference on computer vision and
520 pattern recognition. pp. 7122–7131 (2018) 3
- 521 32. Kocabas, M., Huang, C.H.P., Hilliges, O., Black, M.J.: Pare: Part attention regres-
522 sor for 3d human body estimation. In: Proceedings of the IEEE/CVF International
523 Conference on Computer Vision. pp. 11127–11137 (2021) 3
- 524 33. Kocabas, M., Huang, C.H.P., Tesch, J., Müller, L., Hilliges, O., Black, M.J.: Spec:
525 Seeing people in the wild with an estimated camera. In: Proceedings of the
526 IEEE/CVF International Conference on Computer Vision. pp. 11035–11045
527 (2021) 3
- 528 34. Li, Y., Zhang, S., Wang, Z., Yang, S., Yang, W., Xia, S.T., Zhou, E.: Token-
529 pose: Learning keypoint tokens for human pose estimation. In: Proceedings of the
530 IEEE/CVF International conference on computer vision. pp. 11313–11322 (2021)
531 3
- 532 35. Lin, J., Lee, G.H.: Multi-view multi-person 3d pose estimation with plane sweep
533 stereo. In: Proceedings of the IEEE/CVF Conference on Computer Vision and
534 Pattern Recognition. pp. 11886–11895 (2021) 3, 4
- 535 36. Liu, Y., Gall, J., Stoll, C., Dai, Q., Seidel, H.P., Theobalt, C.: Markerless motion
536 capture of multiple characters using multiview image segmentation. IEEE trans-
537 actions on pattern analysis and machine intelligence **35**(11), 2720–2735 (2013) 3
- 538 37. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned
539 multi-person linear model. ACM transactions on graphics (TOG) **34**(6), 1–16
540 (2015) 2, 3, 4, 5, 6
- 541 38. Luvizon, D.C., Habermann, M., Golyanik, V., Kortylewski, A., Theobalt, C.: Scene-
542 aware 3d multi-human motion capture from a single camera. In: Computer Graph-
543 ics Forum. vol. 42, pp. 371–383. Wiley Online Library (2023) 3
- 544 39. Martinez, J., Hossain, R., Romero, J., Little, J.J.: A simple yet effective baseline for
545 3d human pose estimation. In: Proceedings of the IEEE international conference
546 on computer vision. pp. 2640–2649 (2017) 3
- 547 40. Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.P., Xu,
548 W., Casas, D., Theobalt, C.: Vnect: Real-time 3d human pose estimation with a
549 single rgb camera. Acm transactions on graphics (tog) **36**(4), 1–14 (2017) 3
- 550 41. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng,
551 R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Commu-
552 nications of the ACM **65**(1), 99–106 (2021) 7
- 553 42. Muller, L., Osman, A.A.A., Tang, S., Huang, C.H.P., Black, M.J.: On self-contact
554 and human pose. In: Proceedings of the IEEE/CVF Conference on Computer Vi-
555 sion and Pattern Recognition (CVPR). pp. 9990–9999 (June 2021) 2
- 556 43. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with
557 a multiresolution hash encoding. ACM Transactions on Graphics (ToG) **41**(4), 1–
558 15 (2022) 2, 5, 6
- 559 44. Osman, A.A., Bolkart, T., Black, M.J.: Star: Sparse trained articulated human
560 body regressor. In: Computer Vision–ECCV 2020: 16th European Conference,
561 Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16. pp. 598–613. Springer
562 (2020) 4

- 563 45. Patel, P., Huang, C.H.P., Tesch, J., Hoffmann, D.T., Tripathi, S., Black, M.J.: 563
 564 Agora: Avatars in geography optimized for regression analysis. In: Proceedings 564
 565 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 565
 566 13468–13478 (2021) 3 566
- 567 46. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A., Tzionas, D., 567
 568 Black, M.J.: Expressive body capture: 3d hands, face, and body from a single im- 568
 569 age. In: Proceedings of the IEEE/CVF conference on computer vision and pattern 569
 570 recognition. pp. 10975–10985 (2019) 14 570
- 571 47. Peng, S., Zhang, Y., Xu, Y., Wang, Q., Shuai, Q., Bao, H., Zhou, X.: Neural 571
 572 body: Implicit neural representations with structured latent codes for novel view 572
 573 synthesis of dynamic humans. In: Proceedings of the IEEE/CVF Conference on 573
 574 Computer Vision and Pattern Recognition. pp. 9054–9063 (2021) 4 574
- 575 48. Pishchulin, L., Insafutdinov, E., Tang, S., Andres, B., Andriluka, M., Gehler, P.V., 575
 576 Schiele, B.: Deepcut: Joint subset partition and labeling for multi person pose 576
 577 estimation. In: Proceedings of the IEEE conference on computer vision and pattern 577
 578 recognition. pp. 4929–4937 (2016) 3 578
- 579 49. Reddy, N.D., Guigues, L., Pishchulin, L., Eledath, J., Narasimhan, S.G.: Tesse- 579
 580 track: End-to-end learnable multi-person articulated 3d pose tracking. In: Proceed- 580
 581 ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 581
 582 pp. 15190–15200 (2021) 4 582
- 583 50. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing 583
 584 hands and bodies together. arXiv preprint arXiv:2201.02610 (2022) 14 584
- 585 51. Saito, S., Yang, J., Ma, Q., Black, M.J.: Scanimate: Weakly supervised learning of 585
 586 skinned clothed avatar networks. In: Proceedings of the IEEE/CVF Conference on 586
 587 Computer Vision and Pattern Recognition. pp. 2886–2897 (2021) 4 587
- 588 52. Shuai, Q., Geng, C., Fang, Q., Peng, S., Shen, W., Zhou, X., Bao, H.: Novel 588
 589 view synthesis of human interactions from sparse multi-view videos. In: ACM SIG- 589
 590 GRAPH 2022 Conference Proceedings. SIGGRAPH '22, Association for Comput- 590
 591 ing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3528233.3530704> 3, 4, 8, 9 591
- 592 53. Shuai, Q., Yu, Z., Zhou, Z., Fan, L., Yang, H., Yang, C., Zhou, X.: Reconstructing 592
 593 close human interactions from multiple views. ACM Transactions on Graphics (dec 593
 594 2023). <https://doi.org/10.1145/3618336> 3, 4 594
- 595 54. Song, J., Chen, X., Hilliges, O.: Human body model fitting by learned gradient de- 595
 596 scent. In: European Conference on Computer Vision. pp. 744–760. Springer (2020) 596
 597 14 597
- 598 55. Sun, Y., Bao, Q., Liu, W., Fu, Y., Black, M.J., Mei, T.: Monocular, one-stage, 598
 599 regression of multiple 3d people. In: Proceedings of the IEEE/CVF international 599
 600 conference on computer vision. pp. 11179–11188 (2021) 3 600
- 601 56. Sun, Y., Liu, W., Bao, Q., Fu, Y., Mei, T., Black, M.J.: Putting people in 601
 603 their place: Monocular regression of 3d people in depth. In: Proceedings of the 603
 604 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13243– 604
 605 13252 (2022) 3 605
- 606 57. Tu, H., Wang, C., Zeng, W.: Voxelpose: Towards multi-camera 3d human pose 606
 607 estimation in wild environment. In: Computer Vision–ECCV 2020: 16th European 607
 608 Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16. pp. 197– 608
 609 212. Springer (2020) 2, 4, 8, 9 609
- 610 58. Wang, C., Li, J., Liu, W., Qian, C., Lu, C.: Hmor: Hierarchical multi-person ordi- 610
 611 nical relations for monocular multi-person 3d pose estimation. In: Computer Vision– 611
 612 ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Pro- 612
 613 ceedings, Part III 16. pp. 242–259. Springer (2020) 3 613

- 614 59. Wang, S., Schwarz, K., Geiger, A., Tang, S.: Arah: Animatable volume rendering
615 of articulated human sdbs. In: European conference on computer vision. pp. 1–19.
616 Springer (2022) 4 614
- 617 60. Wang, T., Zhang, J., Cai, Y., Yan, S., Feng, J.: Direct multi-view multi-person
618 3d pose estimation. In: Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P.,
619 Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems. vol. 34,
620 pp. 13153–13164. Curran Associates, Inc. (2021), https://proceedings.neurips.cc/paper_files/paper/2021/file/6da9003b743b65f4c0cccd295cc484e57-Paper.pdf 2, 3,
621 4, 10, 11 615
- 622 61. Weng, C.Y., Curless, B., Srinivasan, P.P., Barron, J.T., Kemelmacher-Shlizerman,
623 I.: Humannerf: Free-viewpoint rendering of moving people from monocular video.
624 In: Proceedings of the IEEE/CVF conference on computer vision and pattern
625 Recognition. pp. 16210–16220 (2022) 4, 8 616
- 626 62. Wu, S., Jin, S., Liu, W., Bai, L., Qian, C., Liu, D., Ouyang, W.: Graph-based
627 3d multi-person pose estimation using multi-view images. In: Proceedings of the
628 IEEE/CVF international conference on computer vision. pp. 11148–11157 (2021)
629 2, 3, 4, 10, 11 617
- 630 63. Wu, S., Jin, S., Liu, W., Bai, L., Qian, C., Liu, D., Ouyang, W.: Graph-based 3d
631 multi-person pose estimation using multi-view images. In: ICCV (2021) 10 618
- 632 64. Xu, H., Bazavan, E.G., Zanfir, A., Freeman, W.T., Sukthankar, R., Sminchisescu,
633 C.: Ghum & ghuml: Generative 3d human shape and articulated pose models.
634 In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
635 Recognition. pp. 6184–6193 (2020) 4 619
- 636 65. Ye, H., Zhu, W., Wang, C., Wu, R., Wang, Y.: Faster voxelpose: Real-time 3d
637 human pose estimation by orthographic projection. In: European Conference on
638 Computer Vision. pp. 142–159. Springer (2022) 2, 3, 4, 10, 11 620
- 639 66. Yin, Y., Guo, C., Kaufmann, M., Zarate, J.J., Song, J., Hilliges, O.: Hi4d: 4d in-
640 stance segmentation of close human interaction. In: Proceedings of the IEEE/CVF
641 Conference on Computer Vision and Pattern Recognition. pp. 17016–17027 (2023)
642 5, 9, 11 621
- 643 67. Zanfir, A., Marinou, E., Sminchisescu, C.: Monocular 3d pose and shape estima-
644 tion of multiple people in natural scenes-the importance of multiple scene
645 constraints. In: Proceedings of the IEEE Conference on Computer Vision and Pattern
646 Recognition. pp. 2148–2157 (2018) 3 622
- 647 68. Zhang, H., Tian, Y., Zhou, X., Ouyang, W., Liu, Y., Wang, L., Sun, Z.: Pymaf: 3d
648 human pose and shape regression with pyramidal mesh alignment feedback loop.
649 In: Proceedings of the IEEE/CVF International Conference on Computer Vision.
650 pp. 11446–11456 (2021) 3 623
- 651 69. Zhang, J., Liu, X., Ye, X., Zhao, F., Zhang, Y., Wu, M., Zhang, Y., Xu, L., Yu, J.:
652 Editable free-viewpoint video using a layered neural representation. ACM Trans-
653 actions on Graphics (TOG) **40**(4), 1–18 (2021) 3, 4, 6, 7 624
- 654 70. Zhang, Y., Wang, C., Wang, X., Liu, W., Zeng, W.: Voxeltrack: Multi-person 3d
655 human pose estimation and tracking in the wild. IEEE Transactions on Pattern
656 Analysis and Machine Intelligence **45**(2), 2613–2626 (2022) 4 625
- 657 71. Zhang, Y., An, L., Yu, T., Li, X., Li, K., Liu, Y.: 4d association graph for realtime
658 multi-person motion capture using multiple video cameras. In: Proceedings of the
659 IEEE/CVF conference on computer vision and pattern recognition. pp. 1324–1333
660 (2020) 2, 3, 8, 10, 11 626
- 661 72. Zhang, Y., Li, Z., An, L., Li, M., Yu, T., Liu, Y.: Lightweight multi-person total
662 motion capture using sparse multi-view cameras. In: Proceedings of the IEEE/CVF
663 International Conference on Computer Vision. pp. 5560–5569 (2021) 4 627
- 664 664

- 665 73. Zheng, Y., Shao, R., Zhang, Y., Yu, T., Zheng, Z., Dai, Q., Liu, Y.: Deepmulticap: 665
666 Performance capture of multiple characters using sparse multiview cameras. In: 666
667 Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 667
668 6239–6249 (2021) 4, 5 668
- 669 74. Zheng, Z., Zhao, X., Zhang, H., Liu, B., Liu, Y.: Avatarrex: Real-time expressive 669
670 full-body avatars. arXiv preprint arXiv:2305.04789 (2023) 14 670
- 671 75. Zhou, H., Hong, C., Han, Y., Huang, P., Zhuang, Y.: Mh pose: 3d human pose 671
672 estimation based on high-quality heatmap. In: 2021 IEEE International Conference 672
673 on Big Data (Big Data). pp. 3215–3222. IEEE (2021) 4 673
- 674 76. Zhou, Z., Shuai, Q., Wang, Y., Fang, Q., Ji, X., Li, F., Bao, H., Zhou, X.: Quickpose: 674
675 Real-time multi-view multi-person pose estimation in crowded scenes. In: ACM 675
676 SIGGRAPH 2022 Conference Proceedings. pp. 1–9 (2022) 3 676